

Linguistics and Literature Review (LLR)

Volume , Issue 1, March 2015

Journal DOI:

Issue DOI:

ISSN: 2221-6510 (Print) 2409-109X (Online) Journal homepage: <http://journals.umt.edu.pk/llr/Home.aspx>

Towards Sindhi Corpus Construction

Mutee U Rahman

To cite to this article: Mutee U Rahman (2015). Towards Sindhi Corpus Construction, *Linguistics and Literature Review* 1(1): 39- 48.

To link to this article:

Published online: March 31, 2015

Article QR Code:



A publication of the
Department of English Language and Literature
School of Social Sciences and Humanities
University of Management and Technology
Lahore, Pakistan

Towards Sindhi Corpus Construction

Mutee U Rahman

Department of Computer Science, Isra University - Hyderabad, Pakistan

ABSTRACT

The paper discusses the current state of Sindhi corpus construction in detail. Sindhi corpus development issues including corpus acquisition, preprocessing, and tokenization are discussed in detail. Preliminary results and observations which include letter unigram, bigram and trigram frequencies; word frequencies and word bigram frequencies are presented. Current state of Sindhi corpus with its limitations and future work is also discussed. The paper also explores the orthography and script of Sindhi language with reference to corpus development.

Keywords: corpus construction, unigram, bigram, trigram frequencies orthography, script

Introduction

Sindhi is one of the major languages of Pakistan spoken by approximately 30-40 million people (Sindhi Language Authority, 2009), (Collie. J., 2009). Sindhi is being frequently used on internet. Sindhi blogs, literary websites, online newspapers and discussion forums are increasing day by day. After Urdu, Sindhi is the second largest written language of Pakistan. Despite of its online usage and popularity only few language processing resources are available for NLP researchers which include lexicon, fonts and simple word processors. The development of Sindhi language processing resources like linguistic corpora and comprehensive computational lexicon are not even initiated. Sindhi is being written in Persio- Arabic (پښتو), Devnagri (सिन्धी) and roman (sindhi) scripts. Persio- Arabic script is most common script for Sindhi writings in Pakistan and India. Devnagri script is also being used for Sindhi writing in India. Roman script (though not yet standardized) is also getting popularity. Very few written documents are available in roman script but it is being used frequently for communications on internet and cell phones and other smart devices. Due to the fact that most of the online and offline written material of Sindhi is available in Persio-Arabic script Sindhi corpus being constructed is in Persio-Arabic script using UTF-16 encoding.

Following sections discuss the existing work in Pakistani language corpora, orthography and script of Sindhi Language, corpus construction issues, corpus acquisition, pre-processing, tokenization and results of preliminary statistical analysis. Finally the future work is discussed along-with conclusion.

Previous work

Apart from fonts, keyboard design (Bhurgri, 2010a) and few digital dictionaries (CRULP, 2010a) Sindhi language processing resources are not available publically. Studies or development projects for resources like linguistic corpora and comprehensive computational lexicon are not even initiated. Various research organizations and individuals are working for the development of linguistic corpora of different Pakistani languages. For Urdu EMILLE (Mentery et al., 2000), Baker Riaz corpus (Becker & Riaz, 2002), jang newspaper corpus (Hussain, 2008), and parallel English Urdu and Nepali corpus (CRULP, 2010b) are some key examples. For Pashto the projects include BBN Byblos Pashto OCR System (Decerbo et al., 2004) and Machine readable Pashto text corpus being developed at University of Peshawar (Khan & Zuhra, 2007). The first Punjabi language corpus was developed by Central Institute of Indian Languages (CIIL) India (Lehal, 2009). Hindi and Punjabi parallel corpus developed by CDAC Noida is another useful linguistic corpora available. One cannot find such type of linguistic corpora for Sindhi, Balouchi, Siraiki and many other Pakistani languages. In contrast to other Pakistani languages (Excluding Urdu) Sindhi text in electronic format is easily available and is being continuously collected for corpus under discussion.

Orthography and script of Sindhi language

Sindhi is written in Persio-Arabic script based on extended Arabic character set in Naskh style. Sindhi alphabet is comprised of 52 letters shown in figure 1. The alphabet contains basic letters like پ, ٺ, ڏ, ڙ and secondary letters like ڄ, ڳ and ڇ which are aspirated versions of چ and گ.

چ	ج	ٺ	ٺ	ڌ	ت	پ	پ	ب	ا
ʃ	dʒ	tʰ	t	tʰ	t	bʰ	ʈ	b	
ڌ	د	خ	ح	ڄ	ڄ	پ	ٺ	ڇ	ڇ
dʰ	d	x	h	tʃʰ	tʃ	p	s	ɳ	dʒʰ
ص	ش	س	ز	ڙ	ر	ذ	ڍ	ڍ	ڏ
s	ʃ	s	z	ʈ	r	z	dʱ	d	dʰ
ڪ	ڪ	ق	ڦ	ف	غ	ع	ظ	ط	ض
kʰ	k	k	pʰ	f	ɣ	ʕ	z	t	z
ء	ھ	و	ڻ	ن	م	ل	گ	ڳ	گ
	h		ɳ	n	m	l	ŋ	gʰ	g
									ي

Figure 1. Sindhi alphabet

Sindhi words always end in a vowel (Rahman, 2009); this vocalic ending is optionally marked by diacritics in written text. Diacritics are also used inside words to represent additional vocal features. Absence of diacritics in written text sometimes causes semantic ambiguities. For instance the word ڀڄڻ (to push) and ڀڄڻ (bog) are semantically ambiguous without diacritics. Diacritics used in Sindhi are shown in Figure 2.

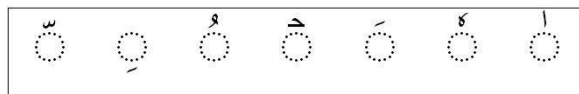


Figure 2. Diacritics used in Sindhi.

Sindhi has its own numerals based on Persio-Arabic numerals shown in figure 3. Use of Hindu-Arabic numerals is also very common in Sindhi writings. Special symbols shown in figure 3 are also used in Sindhi written text.

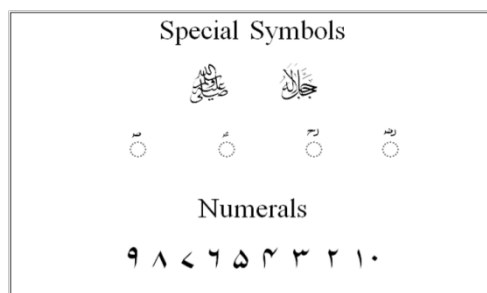


Figure 3. Special symbols and numerals used in Sindhi written text.

Sindhi corpus development

After Unicode support and Unicode based Sindhi keyboard design (Bhurgri, 2010b) availability of Unicode based Sindhi text on Internet is increasing day by day. Key factor behind the motivation of Sindhi corpus construction is availability of online text in Sindhi newspapers, blogs, and literary websites and discussion forums. Despite of the fact that available online resources do not provide huge amount of text but they are increasing day by day and corpus is being collected continuously. Software routines for preprocessing, normalization, tokenization and frequency calculation are implemented in C# using Microsoft.net framework libraries.

Corpus acquisition

Data is gathered from various domains which include news, blogs, literature, essays, and letters. Different subdomains include current affairs, sports, showbiz, short stories, discussions and opinions. Sources of data collection are shown in Table 1.

Table 1. Sources of data collection

Source	URL(s)
Daily Kawish	http://www.thekawish.com
Daily Awami Awaz	http://www.awamiawaz.com
Daily Ibrat	http://dailyibrat.com
Blogs	http://shikarpuri.wordpress.com
Literary Writings	http://voiceofsindh.net http://sindhshamat.com

Preprocessing and normalization

Almost all data gathered was already in Unicode format but nevertheless all the collected text is converted into standard UTF-16 encoding. Letters represented by multiple Unicode points and equivalent representations of composed and decomposed form (Hussain & Durrani, 2008) are reduced to same underlying form. Letters with aspirated versions like گَ which are combinations of two Unicode characters (for instance گ and َ in case of گَ) are considered single letters while dealing with text processing.

Tokenization

For tokenization white spaces, punctuation markers, special symbols (like \$, %, # etc.) and digits are used as word boundaries. White space word boundary consideration caused problem of embedded space word breaking (For example the single word سڌت ڊيٽ سڌت is divided into two words سڌت and ڊيٽ) is tackled out by using the same technique used for Urdu (Ijaz & Hussain, 2007). Another problem in Sindhi word tokenization occurs when two special words ۽ (in) and ۽ (and) occurred without space like ۽ ۽ ۽ (me: mila:i a) and this was tokenized as a single word. Also in case of ۽ ۽ ۽ (kita:baainqalama (book and pen) in which three words without space are there and were tokenized as single word. Same problem was observed with all the words with non-connective ending like ۽ ۽ ۽ (hi:rap:i:a (drink milk) or starting letters سڌت ۽ ۽ (sindh'aander(in Sindh). Semiautomatic (software based + manual) approach was used to overcome this problem.

Results and observations

A total of 4.1 million word corpus analyzed quantitatively. This preliminary analysis includes letter frequency analysis, letter bigram analysis, letter trigram analysis, word frequency

analysis, and word bigram analysis. These quantitative results are discussed in following sections.

Letter frequencies

A total of 13,968,112 characters in the corpus were analyzed while calculating letter frequencies. Along with 52 letters of Sindhi alphabet ^۱was also considered as a single letter because of its use in Sindhi keyboard as a single letter and single Unicode representation. It was observed that most frequently occurred letter was vowel ِ while least frequently occurred letter was consonant گ. Table 2 shows top 20 most frequently occurred letters in Sindhi corpus with their percentage. While analyzing frequencies it was observed that frequency distribution of individual letters in single file of 50,000 or more words was identical to the letter frequency distribution of whole corpus. This similarity can be seen in graphs of figure 4 and 5.

Letter bigram and trigram frequencies were also analyzed. It can be seen that almost 50% of top 20 most frequent bigrams are valid two letter words like ۱ا، ۱ج، ۱ي and ۱ک. Same is the case with trigrams where this ratio is more than 60%. Top 20 most frequent bigram and trigram percentages are shown in Tables 3 and 4 respectively.

Table 2. Top 20 most frequent letters.

S.No.	letter	Percent	S.No.	Letter	Percent
1	ي	13.77%	11	ڪ	3.25%
2	ا	11.42%	12	ط	3.23%
3	ى	8.99%	13	د	2.50%
4		7.84%	14	ة	2.00%
5		6.27%	15	پ	1.80%
6	س	6.15%	16	آ	1.18%
7	م	3.73%	17	ڻ	1.16%
8	ج	3.64%	18	ک	1.16%
9	ل	3.30%	19	ع	0.99%
10	ت	3.26%	20	ٺ	0.94%

Table 3. Top 20 bigrams in Sindhi corpus

S.No.	Bigram	Percent	S.No.	Bigram	Percent
1	نا	3.16%	11	ئ	1.18%
2	۱ج	2.55%	12	بی	1.10%
3	یس	1.95%	13	آ	1.10%
4	۱ي	1.80%	14	ج	1.07%
5	ى	1.79%	15	أ	1.02%

Table 4. Top 20 letter trigrams in Sindhi Corpus

S.No.	Trigram	Percent	S.No.	Trigram	Percent
1	يَا	1.40%	11	جُؤ	0.45%
2	يُؤ	1.34%	12	يَا	0.44%
3	يسا	0.81%	13	يڪ	0.44%
4	يئ	0.74%	14	اُ	0.42%
5	يشڪ	0.71%	15	ذُ	0.41%
6	يڪ	0.61%	16	بڻا	0.40%
7	ذئ	0.60%	17	يچُ	0.36%
8	يذُ	0.53%	18	يڏُ	0.35%
9	ساُ	0.47%	19	پُ	0.35%
10	ببه	0.46%	20	ساد	0.35%

Table 5. Top 20 most frequent words in Sindhi corpus

S.No.	word	Percent	S.No.	word	Percent
1	يچ	3.71%	11	يشڪ	0.69%
2	ڄ	2.44%	12	سبع	0.69%
3	ءِ	2.17%	13	نا	0.67%
4	ت	1.78%	14	يڪ	0.63%
5	پا	1.61%	15	بڻ	0.57%
6	پڪ	1.61%	16	پا	0.55%
7	ج	1.50%	17	ءِلا	0.51%
8	بڻ	1.05%	18	يُ	0.50%
9	ث	0.82%	19	ؤؤ	0.50%
10	ؤؤ	0.71%	20	يڪ	0.46%

Table 6. Top 10 most frequent word bigrams

S.No.	Word bigram	Percentage
1	ي چ ت	7.52
2	پا ت	6.75
3	يُ يچ	2.66
4	ش ڀڙوڻا تڀ	1.93
5	دع يچ	1.84
6	نا يچ	1.72
7	پڏج ت	1.60
8	يُ يچ	1.60

9	ی پک	1.44
1	ئ یا	1.21
0		

Future work

Corpus is being continuously collected and results are being updated. Currently corpus is simply a UTF16 encoded text collection. Studies are in progress for proper annotations, POS tagging, corpus based lexicon development and n-gram based text categorization.

Sindhi tokenization algorithm need to be worked out for the problems discussed in section 4.3. Due to absence of standard sentence termination punctuation marker in Sindhi; full stop comma and other punctuation markers are used as sentence terminators in Sindhi text writings. Sentence segmentation is another key area to be worked out. More specific Sindhi computational linguistic studies are needed for further development and maturity of corpus. For example currently there is no comprehensive POS tagging algorithm available for Sindhi. Presently available POS tagging algorithm for Sindhi (Mahar & Memon, 2010) need to be analyzed and extended further. Sindhi tag set needs to be designed before POS tagging of the corpus. Qualitative, quantitative improvements, proper annotations and comprehensive statistical analysis are areas to be extensively worked out.

Conclusion

In absence of language processing resources of Sindhi language Sindhi corpus construction project is a valuable initiative. Regardless of its size and preliminary results the corpus in its current state will provide basis for further natural language processing studies of Sindhi language. Letter frequencies including bigram and trigram frequencies provide basis for intelligent text processing and compact keyboard design for cell phones and other smart devices. Word level unigram and bigram frequencies provide basis for spelling corrections and automatic sentence completion applications. Further developments in corpus will be useful for advanced language processing tasks like morphological analysis, syntax analysis, semantic analysis, information retrieval and extraction and machine translation.

References

- Sindhi Language Authority. 2009. *Sindhi Language 2010*. Retrieved from.
<http://www.sindhila.org>
- Becker D., and Riaz K. 2002. A Study in Urdu Corpus Construction. *In the Proceedings of 3rd Workshop on Asian Language Resources and International Standardization at the 19th International Conference on Computational Linguistics* 12: 1-5.
- Bhurgri A. M. 2010, July 7. A Breakthrough in use of Sindhi on Internet *Indus Asia Online Journal* <http://iaoj.wordpress.com/2010/07/07/a-breakthrough-in-use-of-sindhi-on-internet/>

- Bhurgri. A. M. 2010. *Sindhi Web-Keyboard2010*. <http://www.bhurgri.com/>
- Collie. J. 2006. The Sindhi Language. In K. Brown (ed). *Encyclopedia of Language and Linguistics, 2nd Edition* 11: 384-386. Oxford: Elsevier.
- CRULP. 2010. *Sindhi English Dictionary 2010*. <http://www.crulp.org/sed/>
- CRULP Parallel Corpus. 2010. *Urdu, Nepali and English Parallel Corpus 2010*. http://crulp.org/software/ling_resources/UrduNepaliEnglishP-parallelCorpus.htm
- Decerbo, M., MacRostie, E., and Natarajan, P. 2004. The BBN Byblos Pashto OCR system. In *Proceedings of the 1st ACM Workshop on Hardcopy Document Processing*, 29-32. ACM New York: NY, USA, Washington, DC, USA.
- Hussain, S. 2008. Resources for Urdu Language Processing. *The Proceedings of the 6th Workshop on Asian Language Resources, IJCNLP'08*, IIIT Hyderabad, India.
- Hussain. S. and Durrani N. 2008. *Sindhi*. PAN Localization A Study on Collation of Languages from Developing Asia. PAN Localization Project. International Development Research Center Canada.
- Ijaz, M. and Hussain, S. 2007. Corpus Based Urdu Lexicon Development. In the *Proceedings of Conference on Language Technology 07*, University of Peshawar, Peshawar, Pakistan.
- Khan M. A., and Zuhra F.T. 2007. A General-Purpose Monitor Corpus of Written Pashto. In the proceedings of *Conference on Corpus Linguistics*, Birmingham.
- Lehal. G. S. 2009. A Survey of the State of the Art in Punjabi Language Processing. *Language In India* 9(10): 9-23.
- Mahar J.A. and Memon G.Q. 2010. Rule Based Part of Speech Tagging of Sindhi Language, ICSAP, *International Conference on Signal Acquisition and Processing*, 101-106.
- McEnery A.M., and Baker P., Gaizauskas R. and Cunningham H. EMILLE: Building a Corpus of South Asian Languages. *Vivek, A Quarterly in Artificial Intelligence* 13(3): 23-32.
- Rahman. M. 2009. Sindhi Morphology and Noun Inflections. In the proceedings of *Conference on Language and Technology CLT09*, Crulp Lahore, 74-81

